

- today:
- finish variational inference
 - Bayesian
 - model selection & causality

mean field approximation

$$\min_{q \in Q_{MF}} KL(q \parallel p)$$

$$\downarrow$$

$$\{q: q(x) = \prod_i q_i(x_i)\}$$

• $KL(\cdot \parallel p)$ is a convex fct. of q
 but Q_{MF} is non-convex constraint set

e.g. Ising model $M_{ij} = M_i M_j$

[see lecture 22 in 2017 for "marginal polytope" perspective]

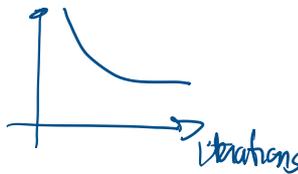


non-convex constraint (on moment parameterization)

[lecture 22 Fall 2017 link](#)

but can monitor progress

$$KL(q^{(t)} \parallel p) + cst.$$



pros & cons of variational methods

vs. Sampling

- ⊕ optimization based
 ⇒ often faster to run & easier to debug

- ⊖ noisy ⇒ harder to debug
- mixing problem for chains

- ⊖ biased estimate

$$\mathbb{E}_{q^{(t)}}[f(z)] \neq \mathbb{E}_p[f(z)]$$

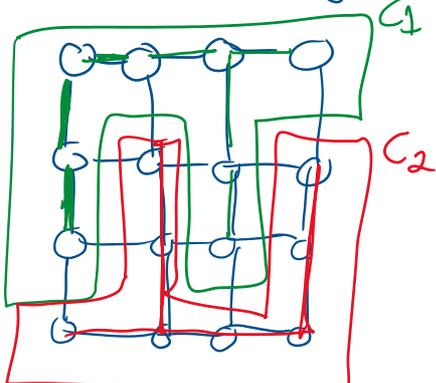
- ⊕ unbiased estimate

$$\mathbb{E}[\mathbb{E}_{q^{(t)}}[f(z)]] = \mathbb{E}_p[f(z)]$$

↳ to make sure chain as mixed with respect to random samples

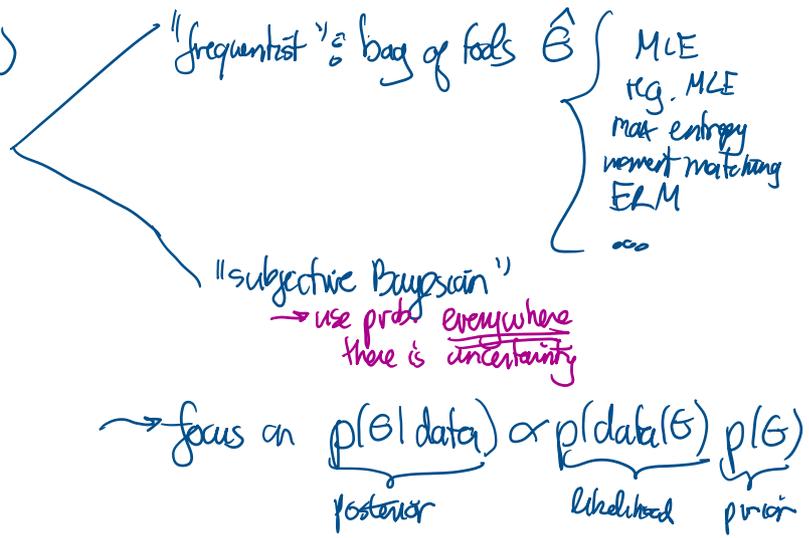
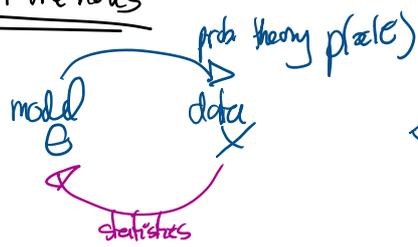
structured mean field

idea $q(z) = \prod_{j=1}^k q_j(z_{C_j})$ where C_1, \dots, C_k is a partition of V and q_j 's are tractable distribution (for example free UGM) for q_j



$$q = q_1 \cdot q_2$$

Bayesian methods



Caricature: • Bayesian is "optimist"
 they think you can get "good" models
 ⇒ obtain a method by doing inference on model

hyperparameters for prior → α_0, β_0

• frequentist is "pessimist" → use analysis tools

Example: biased coin

$X_i | \theta \sim \text{Bernoulli}(\theta)$



e.g. $\theta \sim \text{Unif}[0,1] = \text{Beta}(1,1)$

$p(\theta) = \text{Beta}(\theta | \alpha_0, \beta_0)$

$p(x_i | \theta) = \theta^{x_i} (1-\theta)^{1-x_i}$

posterior $\propto \left(\prod_i p(x_i | \theta) \right) p(\theta)$
 $= \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} \theta^{\alpha_0 - 1} (1-\theta)^{\beta_0 - 1} \mathbb{1}_{[0,1]}(\theta)$

⇒ $p(\theta | \text{data}) = \text{Beta}(\theta | \alpha_0 + n_1, \beta_0 + n - n_1)$

↳ "conjugate prior" to the Bernoulli likelihood model

Conjugate priors

consider a family F of dist. on θ $F = \{ p(\theta | \alpha) : \alpha \in \Omega \}$

say that F is a "conjugate family" to observation model $p(x | \theta)$

if posterior $p(\theta | x, \alpha) \in F$ for any $x \sim X | \theta$

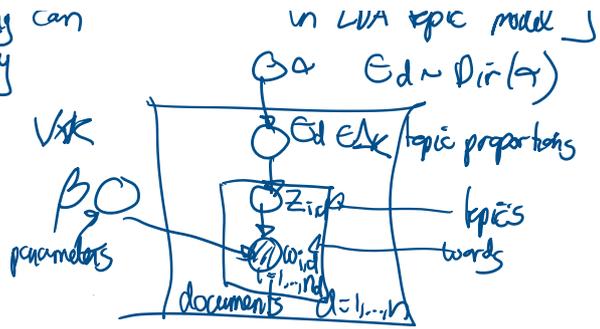
ie. \exists some $\alpha'(x, \alpha)$ s.t. $p(\theta | x, \alpha) = p(\theta | \alpha')$

side note: if use conjugate priors in a DGM then Gibbs sampling can be easy

[eg. this is case in LDA topic model]
 $\rho \sim \text{Dir}(\alpha)$

example here 1 Dirichlet prior is conjugate for multinomial likelihood model

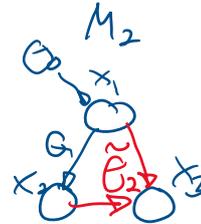
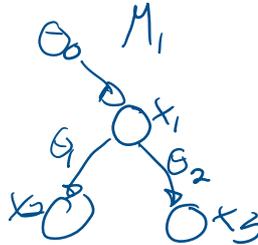
then Gibbs sampling can be easy



15h31

model selection

say we want to choose between 2 DGM



$p(x_3 | x_1; \theta_2)$

$p(x_3 | x_1, x_2; \tilde{\theta}_2)$

[note here that "M1 subset M2"]

"causal notation"

as a frequentist: $\hat{\theta}_{M_1}^{MLE} = \text{argmax}_{\theta_0, \theta_1, \theta_2} \log p(\text{data} | \theta_0, \theta_1, \theta_2, \text{"model"} = M_1)$

$\hat{\theta}_{M_2}^{MLE} = \text{argmax}_{\theta_0, \theta_1, \tilde{\theta}_2} \log p(\text{data} | \theta_0, \theta_1, \tilde{\theta}_2, \text{"model"} = M_2)$
 ↳ different space than θ_2

how to choose between models?

can't compare $\log p(\text{data} | \hat{\theta}_{M_1}^{MLE}, M=M_1)$ vs $\log p(\text{data} | \hat{\theta}_{M_2}^{MLE}, M=M_2)$
 because LHS \leq RHS since $M_1 \subseteq M_2$
 (ie. you would always choose "bigger model")

→ as a frequentist, use cross-validation or validation set
 ie. $\log p(\text{test data} | \hat{\theta}_{M_i}^{MLE}(\text{train data}))$
 $M=M_i$

Bayesian alternatives

free Bayesian ⇒ sum over models (integrate out uncertainty about M_i)

introduce a prior over models $p(M)$

$p(x_{\text{new}} | \mathcal{D}) = \sum_M p(x_{\text{new}} | \mathcal{D}, M) p(M | \mathcal{D})$
 $= \sum_M \left[\int_{\theta \in \Theta_M} p(x_{\text{new}} | \theta, M) p(\theta | \mathcal{D}, M) d\theta \right] p(M | \mathcal{D})$
 (sum over posterior over models)

$$\sum_M \left(\int_{\Theta \in \Theta_M} p(\theta_{new} | \Theta_M, M) p(\theta | D, M) d\theta \right) p(M | D)$$

standard Bayesian predictive dist. for one model (pointing to the integral)
 posterior on θ given data D & model M (pointing to $p(\theta | D, M)$)
 doing model averaging (pointing to the sum)

⊗ in model selection, forced to pick model

⇒ pick model that maximizes $p(M | \text{data})$ or $p(\text{data} | M) p(M)$

$$p(\text{data} | M) = \text{"marginal likelihood"}$$

$$\int_{\Theta} p(\text{data} | \Theta, M) p(\Theta | M) d\Theta$$

to compare two models, look at:

$$\frac{p(M=M_1 | D)}{p(M=M_2 | D)} = \frac{p(D | M_1) p(M_1)}{p(D | M_2) p(M_2)}$$

Bayes factor (pointing to the ratio)
 prior odds (pointing to $p(M_1)/p(M_2)$)

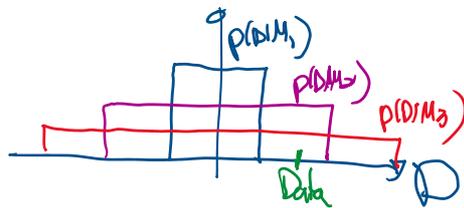
"uniform prior over models" ⇒ then pick among k models M_1, \dots, M_k

by max $p(\text{data} | M = M_i)$

"empirical Bayes" "type II ML"

when # of models is "small", then this approach is "fine" (i.e. won't overfit)

Zoubin's cartoon: suppose $M_1 \subset M_2 \subset M_3$



$p(D|M)$ is normalized over D

vs. $p(D | \hat{\Theta}_{MLE}(D), M)$ [can overfit badly]

but type II ML can still overfit if have too many models

say e.g. $p(D|M) = \delta(D, M)$



how to compute marginal likelihood:

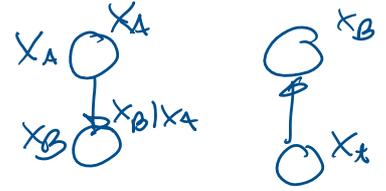
use approximations \leftarrow variational inference / sampling

use approximations \leftarrow variational inference
 sampling

simple approximation \rightarrow Bayesian information criterion (BIC)

Causality:

Structural causal model: graph model + intervention model

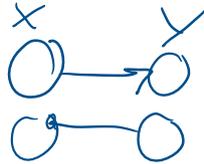
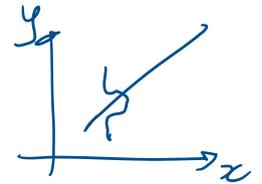


$$p(x) = \prod_{i=1}^n p(x_i | \pi_i, \theta_i)$$

semantic of intervening on node J

$$p(x | \text{intervention on } J) = \left(\prod_{i \neq J} p(x_i | \pi_i, \theta_i) \right) p(x_J | \text{intervention})$$

identify causal direction \leftarrow via parametric assumptions
 via interventions



see thoughts of Bernhard Schölkopf on causality:

<https://arxiv.org/abs/1911.10500>

(and references therein, e.g. his book:)

Elements of Causal Inference, 2017

By Jonas Peters, Dominik Janzing and Bernhard Schölkopf

<https://mitpress.mit.edu/books/elements-causal-inference>

(available for free online)